

Exponential Families

Neil Girdhar

August 4, 2022

1 Exponential families

The *exponential families* are an important class of probability distributions that include the normal, gamma, beta, exponential, Poisson, binomial, and Bernoulli distributions. In this section, based on presentations by Nielsen and Garcia (2011) and Shao (2003, p. 66), we describe some of the exponential families' many exciting properties.

1.1 Definition

For simplicity, we will restrict ourselves to discrete and continuous distributions; the general, measure-theoretic definition (Shao 2003, p. 66) is analogous. A *natural exponential family* is a family of probability distributions parametrized by $\boldsymbol{\eta}$ and whose probability mass function or density function can be decomposed as:

$$f(\mathbf{x} \mid \boldsymbol{\eta}) = \exp(T(\mathbf{x})^T \boldsymbol{\eta} - g(\boldsymbol{\eta}) + h(\mathbf{x})), \quad \mathbf{x} \in \Omega \quad (1)$$

where

- Ω is the *support*,
- $\boldsymbol{\eta}$ are the *natural parameters*,
- $T(\mathbf{x})$ is the *sufficient statistic*,
- $g(\boldsymbol{\eta})$ is the *log-normalizer*, and
- $h(\mathbf{x})$ is the *carrier measure*.

(See Appendix 2 for examples.)

1.1.1 Natural parameters

The decomposition of an exponential family in equation 1 is not unique. Any transformation

$$\boldsymbol{\eta}' = D\boldsymbol{\eta} \quad T' = [D^T]^{-1} T \quad (2)$$

where D is a nonsingular matrix (a bijective linear map) gives another representation of the same natural exponential family.

If $\boldsymbol{\eta}$ were replaced by an arbitrary function $\boldsymbol{\eta}(\boldsymbol{\theta})$ of parameters $\boldsymbol{\theta}$, then the family of probability distributions is called an *exponential family*. We avoid this general form, preferring the so-called *canonical form*.

1.1.2 Sufficient statistic

The sufficient statistic is a vector-valued function of only the outcome \mathbf{x} . Its name is justified by its connection to the maximum entropy formulation (§1.1.5) and to maximum likelihood estimation (§1.3.2).

1.1.3 Log-normalizer

The log-normalizer is a scalar-valued function of only the natural parameters $\boldsymbol{\eta}$. It is so-named because

$$1 = \int_{\mathbf{x}} f(\mathbf{x} | \boldsymbol{\eta}) d\mathbf{x} = \int_{\mathbf{x}} \exp(T(\mathbf{x})^\top \boldsymbol{\eta} - g(\boldsymbol{\eta}) + h(\mathbf{x})) d\mathbf{x} \quad (3)$$

\Downarrow

$$g(\boldsymbol{\eta}) = \log \int_{\mathbf{x}} \exp(T(\mathbf{x})^\top \boldsymbol{\eta} + h(\mathbf{x})) d\mathbf{x}. \quad (4)$$

The log-normalizer is strictly convex and smooth (infinitely differentiable) (Nielsen and Nock 2011).

1.1.4 Carrier measure

The carrier measure is a scalar-valued function of only the outcome \mathbf{x} . In the measure-theoretic presentation of exponential families (Shao 2003, p. 66), the carrier measure truly is a *measure* on the support. The measure-theoretic intuition is analogous to Shannon’s description of continuous entropy (§??): the carrier measure is an assumed standard that weights each small volume of the domain by $\exp(h(\mathbf{x}))$. It represents prior knowledge about the parametrization of the support.

Many formulae are simplified when the carrier measure is zero, in which case it is called a *standard carrier measure*. For continuous distributions, this can always be achieved by a change of variables; for discrete distributions, the carrier measure is rarely a standard carrier measure, and nothing can be done to make it so.

1.1.5 Maximum entropy formulation

The exponential families are motivated in such situations: Suppose that we make independent realizations of a random variable $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, but we only know (1) the expected sufficient statistic $\mathbb{E}(T(\mathbf{x}))$ for some known function T , (2) the support of the realizations Ω , and (3) optionally some prior carrier measure h on the space. Then, Jaynes (1957) avers the *principle of maximum entropy*: “the maximum-entropy distribution may be asserted for the positive reason that it is uniquely determined as the one which is maximally noncommittal with regard to missing information, instead of the negative one that there was no reason to think otherwise. Thus the concept of entropy supplies the missing criterion of choice...” Gokhale (1975) showed that these constraints uniquely lead to a maximum-entropy exponential family with the given sufficient statistic and support.

For example, given a mean μ , and a variance σ^2 , the support of the reals, and assuming standard carrier measure, one is spared the maximum-entropy calculation and can arrive directly at the exponential family distribution function:

$$f(x) \propto \exp \left(\left[\begin{array}{c} x \\ (x - \mu)^2 \end{array} \right]^\top \boldsymbol{\eta} \right) \quad (5)$$

Normalizing this function leads to the normal distribution’s density function (§2.1.1). The parameters $\boldsymbol{\eta}$ are uniquely determined by the given mean and variance.

1.2 The natural parametrization

The *natural parametrization* of an exponential family is the vector space for combining and scaling evidence from independent sources. The natural parametrization is the one that specifies elements of the exponential family using natural parameters (§1.1.1).

1.2.1 Bayesian evidence combination

For example, consider that a friend flips a coin four times and secretly records the result. His belief over the coin's bias is distributed X_1 with natural parameters $\boldsymbol{\eta}_{X_1}$. Then, you flip the coin once and record the result yielding a belief distributed X_2 with natural parameters $\boldsymbol{\eta}_{X_2}$. Given the coin, your beliefs are independent. The “Bayesian evidence combination” operation (Figure 1) aggregates such independent information by summing the natural parameters. This is because the combined belief

$$f(\mathbf{x} \mid \boldsymbol{\eta}_{X_1}, \boldsymbol{\eta}_{X_2}) \propto \exp(T(\mathbf{x})^\top \boldsymbol{\eta}_{X_1} - g(\boldsymbol{\eta}) + h(\mathbf{x})) \exp(T(\mathbf{x})^\top \boldsymbol{\eta}_{X_2} - g(\boldsymbol{\eta})) \quad (\text{by equation 1}). \quad (6)$$

(The decomposition into a product is by independence, and we take care not to double-count the carrier measure h .)

$$= \exp(T(\mathbf{x})^\top (\boldsymbol{\eta}_{X_1} + \boldsymbol{\eta}_{X_2}) - g(\boldsymbol{\eta}) + h(\mathbf{x})). \quad (7)$$

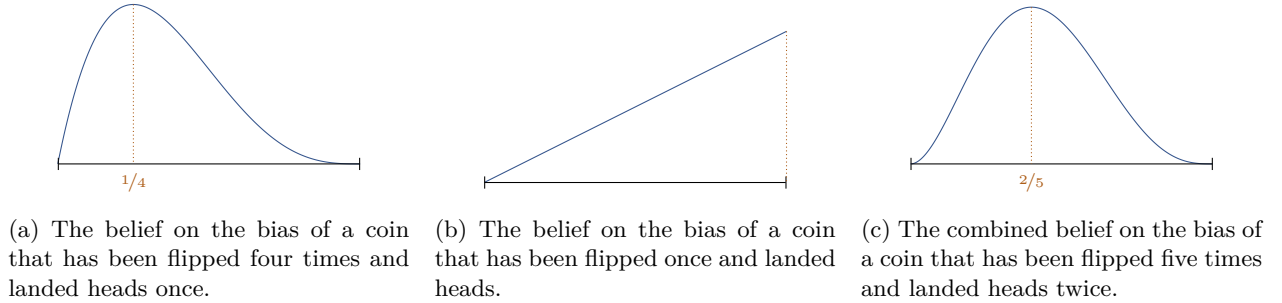


Figure 1: Bayesian evidence combination with beta-distributed (§2.5.1) beliefs over the bias of a coin.

1.2.2 Bayesian evidence scaling

If you value the opinion from a friend more than your own, it is as if n friends provided identical, but independent information. Reasoning from §1.2.1, “Bayesian evidence scaling” (Figure 2) corresponds to scaling in the space of natural parameters.

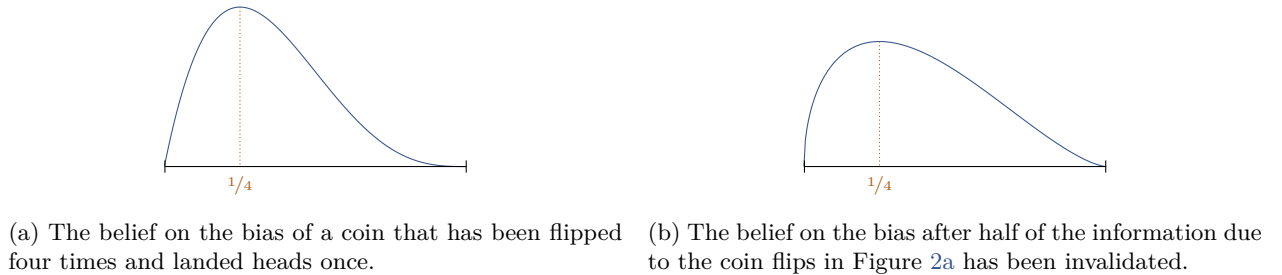


Figure 2: Bayesian evidence scaling with beta-distributed (§2.5.1) beliefs over the bias of a coin.

1.2.3 Bayesian evidence combination is better than product of experts

Hinton (2002) calls a similar operation—the pointwise product of probability measures—a *product of experts*. For an exponential family, this operation is equivalent to “Bayesian evidence combination” except when the carrier measure (§1.1.4) is nonzero. In that case, the carrier measure, which represents prior knowledge about the parametrization of the support, is double-counted. In other words, Bayesian evidence combination is invariant under reparametrization unlike *product of experts*.

1.3 The expectation parametrization

Suppose we have a random variable X distributed according to a distribution in family \mathcal{F} (which is a natural exponential family). Then, X has an *expectation parametrization*, which is the one whose parameters are the expected sufficient statistic

$$\boldsymbol{\chi} \triangleq \mathbb{E}(T(X)). \quad (8)$$

This parametrization is convenient for *parametric density estimation*: the problem of estimating a distribution’s parameters given its realizations.

Like the natural parametrization (equation 2), the expectation parametrization is unique up to a bijective linear map. Unlike the natural parametrization, the expectation parametrization does not have meaningful vector space operations (constant scaling and summation); instead, only weighted averages are meaningful.

1.3.1 Conjugate prior distribution

If the random variable X is distributed according to an exponential family distribution with unknown natural parameters $\boldsymbol{\eta}$, then independent realizations $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ of X induce a likelihood over $\boldsymbol{\eta}$: $\mathcal{L}(\boldsymbol{\eta} \mid \mathbf{x}_1, \dots, \mathbf{x}_n)$. This distribution must belong to a family \mathcal{F}' (called the conjugate prior of \mathcal{F}) that is also an exponential family since

$$\mathcal{L}(\boldsymbol{\eta} \mid \mathbf{x}_1, \dots, \mathbf{x}_n) = f(\mathbf{x}_1, \dots, \mathbf{x}_n \mid \boldsymbol{\eta}) \quad (9)$$

$$\propto \prod_i f(\mathbf{x}_i \mid \boldsymbol{\eta}) \quad (\text{by independence}) \quad (10)$$

$$= \prod_i \exp\left(T(\mathbf{x}_i)^\top \boldsymbol{\eta} - g(\boldsymbol{\eta}) + h(\mathbf{x}_i)\right) \quad (11)$$

(Since \mathbf{x}_i are fixed, $\prod_i \exp(h(\mathbf{x}_i))$ is constant)

$$\propto \prod_i \exp\left(T(\mathbf{x}_i)^\top \boldsymbol{\eta} - g(\boldsymbol{\eta})\right) \quad (12)$$

$$= \exp\left(\left(\sum_i T(\mathbf{x}_i)\right)^\top \boldsymbol{\eta} - ng(\boldsymbol{\eta})\right) \quad (13)$$

$$= \exp(T'(\boldsymbol{\eta})^\top \boldsymbol{\eta}') \quad (14)$$

where

$$T'(\boldsymbol{\eta}) = \begin{bmatrix} \boldsymbol{\eta} \\ g(\boldsymbol{\eta}) \end{bmatrix} \quad \boldsymbol{\eta}' = \begin{bmatrix} \sum_i T(\mathbf{x}_i) \\ -n \end{bmatrix}. \quad (15)$$

In equation 14, we can see that the vector of *hyperparameters* $\boldsymbol{\eta}'$ are natural parameters of the distribution. Thus, “Bayesian evidence combination” and “Bayesian evidence scaling” correspond to vector addition and scaling of these induced parameters, one of which n is a real-valued pseudo-observation count.

1.3.2 Maximum likelihood distribution

Continuing the parametric density estimation problem from §1.3.1, suppose we have an induced likelihood over $\boldsymbol{\eta}$:

$$\mathcal{L}(\boldsymbol{\eta} \mid \mathbf{x}_1, \dots, \mathbf{x}_n) \propto \exp(T'(\boldsymbol{\eta})^\top \boldsymbol{\eta}'). \quad (14 \text{ revisited})$$

Then, the maximum likelihood distribution of X given the realizations $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ is the mode of equation 14. So,

$$0 = \frac{\partial \exp(T'(\boldsymbol{\eta})^\top \boldsymbol{\eta}')}{\partial \boldsymbol{\eta}} \quad (16)$$

$$= \exp(T'(\boldsymbol{\eta})^\top \boldsymbol{\eta}') \frac{\partial \begin{pmatrix} \begin{bmatrix} \eta_1 \\ \vdots \\ \eta_k \\ -g(\boldsymbol{\eta}) \end{bmatrix}^\top \begin{bmatrix} \sum_i T(\mathbf{x}_i)_1 \\ \vdots \\ \sum_i T(\mathbf{x}_i)_k \\ n \end{bmatrix} \end{pmatrix}}{\partial \boldsymbol{\eta}} \quad (17)$$

$$\Downarrow$$

$$\frac{\partial g(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}} = \frac{\sum_i T(\mathbf{x}_i)}{n} \quad (18)$$

$$\Downarrow$$

$$\boldsymbol{\chi} = \frac{\sum_i T(\mathbf{x}_i)}{n} \quad (\text{by equation 29}). \quad (19)$$

Thus, the maximum likelihood distribution has expectation parameters equal to the expected sufficient statistics of the samples. This motivates the expectation parametrization, and the term *sufficient statistic*.

1.3.3 Aggregating maximum likelihood distributions

Suppose that instead of n independent realizations $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ of our exponential family distribution X (as in §1.3.1), we collect realizations m times. The i th collection yields n_i realizations from which we calculate a maximum likelihood distribution X_i having expectation parameters $\boldsymbol{\chi}_i$ (as per §1.3.2). After the collection, we discard the realization, so that all we have are the n_i s and X_i s. How can we combine these into one maximum likelihood distribution over all $\sum_{i=1}^m n_i$ realizations had been collected.

From equations 15 and 19, we can conclude that the natural parameters of the conjugate prior distribution for each i is

$$\boldsymbol{\eta}'_i = \begin{bmatrix} n_i \boldsymbol{\chi}_i \\ n_i \end{bmatrix}. \quad (20)$$

From §1.2.1, we know that we can combine these into one conjugate prior distribution with parameters:

$$\boldsymbol{\eta}_i = \begin{bmatrix} \sum_{i=1}^m n_i \boldsymbol{\chi}_i \\ \sum_{i=1}^m n_i \end{bmatrix}. \quad (21)$$

From equation 19, we can conclude the maximum likelihood distribution given all of the realizations has expectation parameters:

$$\frac{\sum_{i=1}^m n_i \boldsymbol{\chi}_i}{\sum_{i=1}^m n_i}. \quad (22)$$

Therefore, weighted average in the space of expectation parameters represents combining maximum likelihood distributions as if the realizations they were based on had been aggregated.

1.3.4 Duality of parametrizations

If the random variable X has known natural parameters $\boldsymbol{\eta}$, then Nielsen and Nock (2011) show that the function that converts natural parameters to expectation parameters is the gradient of the log-normalizer:

$$\frac{\partial g(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}} = \frac{\partial \log \int_{\mathbf{x}} \exp(T(\mathbf{x})^\top \boldsymbol{\eta} + h(\mathbf{x})) \, d\mathbf{x}}{\partial \boldsymbol{\eta}} \quad (\text{by equation 4}) \quad (23)$$

$$= \frac{\int_{\mathbf{x}} T(\mathbf{x}) \exp(T(\mathbf{x})^\top \boldsymbol{\eta} + h(\mathbf{x})) \, d\mathbf{x}}{\int_{\mathbf{x}} \exp(T(\mathbf{x})^\top \boldsymbol{\eta} + h(\mathbf{x})) \, d\mathbf{x}} \quad (24)$$

$$= \frac{\int_{\mathbf{x}} T(\mathbf{x}) \exp(T(\mathbf{x})^\top \boldsymbol{\eta} + h(\mathbf{x})) \, d\mathbf{x}}{\exp(g(\boldsymbol{\eta}))} \quad (\text{by equation 4}) \quad (25)$$

$$= \int_{\mathbf{x}} T(\mathbf{x}) \exp(T(\mathbf{x})^\top \boldsymbol{\eta} - g(\boldsymbol{\eta}) + h(\mathbf{x})) \, d\mathbf{x} \quad (26)$$

$$= \int_{\mathbf{x}} T(\mathbf{x}) f(\mathbf{x} \mid \boldsymbol{\eta}) \, d\mathbf{x} \quad (\text{by equation 1}) \quad (27)$$

$$= \mathbb{E}(T(X)) \quad (28)$$

$$= \boldsymbol{\chi} \quad (\text{by equation 8}). \quad (29)$$

1.3.5 Higher moments of the sufficient statistic

The higher moments of the sufficient statistic are the higher-order gradients of the log-normalizer:

$$\nabla_{\boldsymbol{\eta}}^n g(\boldsymbol{\eta}) = \nabla_{\boldsymbol{\eta}}^{n-1} \int_{\mathbf{x}} T(\mathbf{x}) \exp(T(\mathbf{x})^\top \boldsymbol{\eta} - g(\boldsymbol{\eta}) + h(\mathbf{x})) \, d\mathbf{x} \quad (\text{by equation 26}) \quad (30)$$

$$= \int_{\mathbf{x}} T(\mathbf{x}) \otimes \left[\otimes_{i=1}^{n-1} \left(T(\mathbf{x}) - \frac{\partial g(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}} \right) \right] \exp(T(\mathbf{x})^\top \boldsymbol{\eta} - g(\boldsymbol{\eta}) + h(\mathbf{x})) \, d\mathbf{x} \quad (31)$$

$$= \int_{\mathbf{x}} T(\mathbf{x}) \otimes \left[\otimes_{i=1}^{n-1} \left(T(\mathbf{x}) - \frac{\partial g(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}} \right) \right] f(\mathbf{x} \mid \boldsymbol{\eta}) \, d\mathbf{x} \quad (\text{by equation 1}) \quad (32)$$

$$= \mathbb{E} \left(T(X) \otimes \left[\otimes_{i=1}^{n-1} \left(T(X) - \frac{\partial g(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}} \right) \right] \right) \quad (33)$$

$$= \mathbb{E} \left(T(X) \otimes \left[\otimes_{i=1}^{n-1} (T(X) - \mathbb{E}(T(X))) \right] \right) \quad (\text{by equation 28}). \quad (34)$$

So, for example, the covariance matrix of the sufficient statistic is the Hessian of the log-normalizer:

$$\frac{\partial^2 g(\boldsymbol{\eta})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}} = \mathbb{E} \left(T(X) \otimes (T(X) - \mathbb{E}(T(X))) \right) \quad (35)$$

$$= \text{Var}(T(X)). \quad (36)$$

In particular, as described by Efron (1978),

$$\frac{\partial \boldsymbol{\chi}}{\partial \boldsymbol{\eta}} = \frac{\partial^2 g(\boldsymbol{\eta})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}} \quad (\text{by equation 29}) \quad (37)$$

$$= \text{Var}(T(X)) \quad (\text{by equation 36}). \quad (38)$$

1.4 Statistics of exponential families

1.4.1 Information theoretic statistics

Consider a data-generating distribution X , and an approximating distribution Y in the same exponential family, having natural parameters $\boldsymbol{\eta}_X$ and $\boldsymbol{\eta}_Y$, and expectation parameters $\boldsymbol{\chi}_X$ and $\boldsymbol{\chi}_Y$. Their cross entropy

(§??) is

$$\mathcal{H}^\times(X; Y) = - \int_{\mathbf{x}} f_X(\mathbf{x}) \log f_Y(\mathbf{x}) \, d\mathbf{x} \quad (\text{by equation ??}) \quad (39)$$

$$= - \int_{\mathbf{x}} f_X(\mathbf{x}) (T(\mathbf{x})^\top \boldsymbol{\eta}_Y - g(\boldsymbol{\eta}_Y) + h(\mathbf{x})) \, d\mathbf{x} \quad (\text{by equation 1}) \quad (40)$$

$$= -\boldsymbol{\chi}_X^\top \boldsymbol{\eta}_Y + g(\boldsymbol{\eta}_Y) - \mathbb{E}_{\mathbf{x} \sim X} (h(\mathbf{x})) \quad (\text{by equation 29}). \quad (41)$$

The entropy (§??) of X is

$$\mathcal{H}(X) = \mathcal{H}^\times(X; X) \quad (\text{by equation ??}) \quad (42)$$

$$= -\boldsymbol{\chi}_X^\top \boldsymbol{\eta}_X + g(\boldsymbol{\eta}_X) - \mathbb{E}_{\mathbf{x} \sim X} (h(\mathbf{x})) \quad (43)$$

and their relative entropy (§??) is

$$\mathcal{H}^{\text{KL}}(X; Y) = \mathcal{H}^\times(X; Y) - \mathcal{H}(X) \quad (\text{by equation ??}) \quad (44)$$

$$= g(\boldsymbol{\eta}_Y) - g(\boldsymbol{\eta}_X) - (\boldsymbol{\eta}_Y - \boldsymbol{\eta}_X)^\top \boldsymbol{\chi}_X. \quad (45)$$

So, for exponential families, the information theoretic statistics are easily calculated from the natural and expectation parameters (Figure 3).

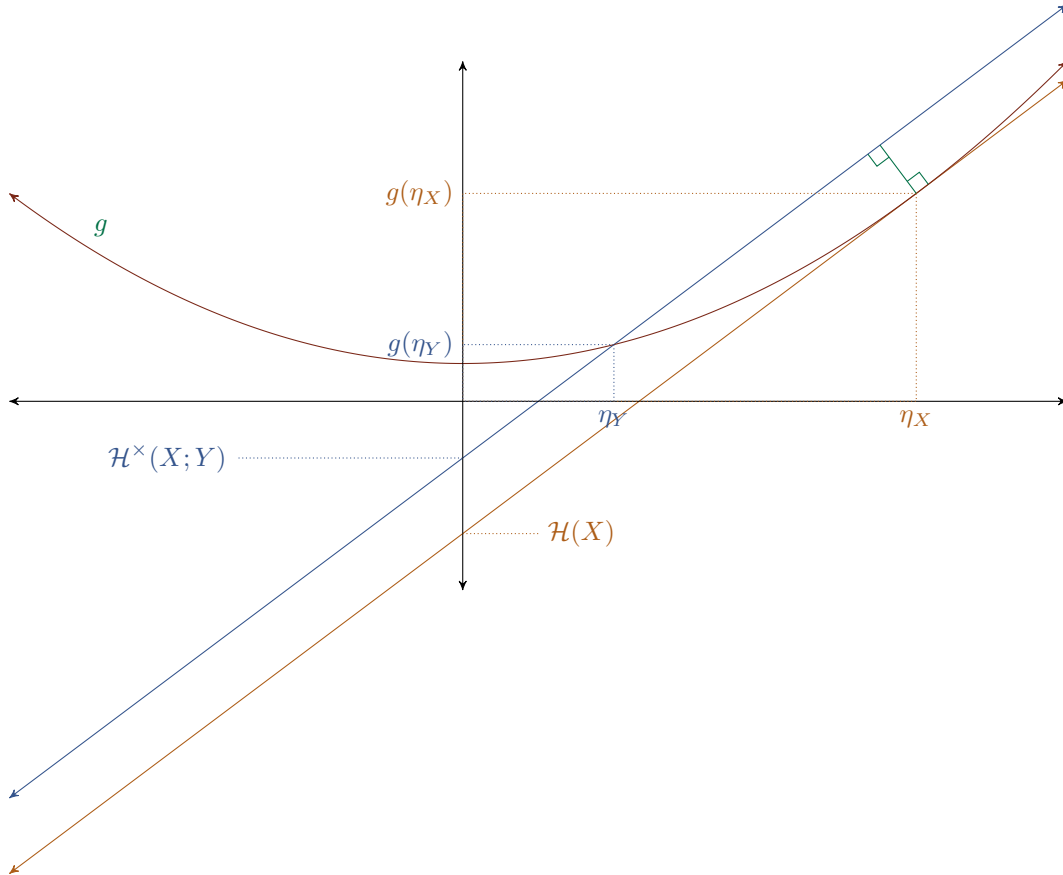


Figure 3: Graphical illustration of the entropy, cross entropy, and relative entropy of exponential families with standard carrier measure adapted from Nielsen and Nock (2011).

1.4.2 Parameter estimation statistics

The statistical score (§??) of $\boldsymbol{\eta}_Y$ given data $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ is

$$\mathcal{V}(\boldsymbol{\eta}_Y \mid \mathbf{x}_1, \dots, \mathbf{x}_n) = \frac{\partial \log \mathcal{L}(\boldsymbol{\eta}_Y \mid \mathbf{x}_1, \dots, \mathbf{x}_n)}{\partial \boldsymbol{\eta}_Y} \quad (46)$$

$$= \frac{\partial \sum_{i=1}^n \log \mathcal{L}(\boldsymbol{\eta}_Y \mid \mathbf{x}_i)}{\partial \boldsymbol{\eta}_Y} \quad (47)$$

$$= \sum_{i=1}^n \frac{\partial (T(\mathbf{x}_i)^\top \boldsymbol{\eta}_Y - g(\boldsymbol{\eta}_Y) + h(\mathbf{x}_i))}{\partial \boldsymbol{\eta}_Y} \quad (\text{by equation 1}) \quad (48)$$

$$= \sum_{i=1}^n \left(T(\mathbf{x}_i) - \frac{\partial g(\boldsymbol{\eta}_Y)}{\partial \boldsymbol{\eta}_Y} \right) \quad (49)$$

$$= \sum_{i=1}^n (T(\mathbf{x}_i) - \boldsymbol{\chi}_Y) \quad (\text{by equation 29}). \quad (50)$$

Therefore, the expected value of the score is

$$\mathbb{E}_{\mathbf{x} \sim X} (\mathcal{V}(\boldsymbol{\eta}_Y \mid \mathbf{x})) = \boldsymbol{\chi}_X - \boldsymbol{\chi}_Y \quad (\text{by equation 8}) \quad (51)$$

$$= - \frac{\partial \mathcal{H}^\times(X; Y)}{\partial \boldsymbol{\eta}_Y} \quad (\text{by equation ??}). \quad (52)$$

The Fisher information (§??) is

$$\mathcal{I}(\boldsymbol{\eta}_Y) = - \mathbb{E}_{\mathbf{x} \sim Y} \left(\frac{\partial^2 \log f(\mathbf{x} \mid \boldsymbol{\eta}_Y)}{\partial \boldsymbol{\eta}_Y \partial \boldsymbol{\eta}_Y} \right) \quad (53)$$

$$= - \mathbb{E}_{\mathbf{x} \sim Y} \left(\frac{\partial^2 (T(\mathbf{x})^\top \boldsymbol{\eta}_Y - g(\boldsymbol{\eta}_Y) + h(\mathbf{x}))}{\partial \boldsymbol{\eta}_Y \partial \boldsymbol{\eta}_Y} \right) \quad (\text{by equation 1}) \quad (54)$$

$$= \mathbb{E}_{\mathbf{x} \sim Y} \left(\frac{\partial^2 g(\boldsymbol{\eta}_Y)}{\partial \boldsymbol{\eta}_Y \partial \boldsymbol{\eta}_Y} \right) \quad (55)$$

$$= \frac{\partial^2 g(\boldsymbol{\eta}_Y)}{\partial \boldsymbol{\eta}_Y \partial \boldsymbol{\eta}_Y}. \quad (56)$$

The Jeffreys prior (§??) for a natural exponential family is thus

$$f(\boldsymbol{\eta}) \propto \sqrt{\det \mathcal{I}(\boldsymbol{\eta})} = \sqrt{\det \frac{\partial^2 g(\boldsymbol{\eta})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}}}. \quad (57)$$

1.5 Altering exponential families

§?? defines *generalized linear models*, which are a kind of *regression* (§??) that makes an exponential family distributional assumption about the targets and uses cross entropy loss. This section explores what happens in the case of three alterations of the assumed exponential family.

Learning in generalized linear models depends only on the gradient of the cross entropy, which is the difference of the expectation parameters of the target values and those of predictions (by equation 51). If we find that the expectation parameters are affected by an alteration, then learning is affected, which means that the model is different. Otherwise, the alteration has no effect on the model.

1.5.1 Transformation of an exponential family

The beta distribution (§2.5.1) is the *conjugate prior* (§1.3.1) of a Bernoulli distribution (§2.2.2) parametrized by a probability $p \in [0, 1]$. If instead we had parametrized the Bernoulli using *odds* $o = \frac{p}{p+1}$, the conjugate prior is *beta-prime*. Therefore, for any beta-distributed X , there is a beta-prime-distributed $\frac{X}{X+1}$. Is regression with a beta distributional assumption the same as regression with a beta-prime assumption?

In general, suppose that we have an exponential family \mathcal{F} with sufficient statistics $T_{\mathcal{F}}$ and carrier measure $h_{\mathcal{F}}$ over support \mathcal{S} . For any distribution $D \in \mathcal{F}$, let $X \sim D$ be a random variable with density f_X and distribution function F_X .

Let $a : \mathcal{S} \rightarrow \mathcal{T}$ be a smooth, invertible function that is independent of the parameters of X and let $Y = a(X)$ be a random variable with density f_Y and distribution function F_Y . The distribution function of Y is

$$F_Y(y) = F_X(a^{-1}(y)) \quad (58)$$

\Downarrow

$$f_Y(y) \triangleq \frac{dF_Y(y)}{dy} = \frac{dF_X(a^{-1}(y))}{dy} \quad (59)$$

$$= f_X(a^{-1}(y)) \frac{da^{-1}(y)}{dy}. \quad (60)$$

Therefore, Y 's distribution belongs to an exponential family \mathcal{G} with sufficient statistics

$$T_{\mathcal{G}}(y) = T_{\mathcal{F}}(a^{-1}(y)), \quad (61)$$

carrier measure

$$h_{\mathcal{G}}(y) = h_{\mathcal{F}}(y) + \log \left(\frac{da^{-1}(y)}{dy} \right), \quad (62)$$

support \mathcal{T} , and the same log-normalizer as \mathcal{F} .

The expectation parameters are unchanged since

$$\mathbb{E}(T_{\mathcal{G}}(Y)) = \mathbb{E}(T_{\mathcal{F}}(a^{-1}(Y))) \quad (\text{by equation 61}) \quad (63)$$

$$= \mathbb{E}(T_{\mathcal{F}}(X)). \quad (64)$$

This shows that changing the distributional assumption of a generalized linear model from \mathcal{F} to \mathcal{G} by smoothly transforming its values has no effect on the model.

1.5.2 Truncation of an exponential family

Linear regression is equivalent to an assumption of normality. However, if the target values are known to be from a subset of the reals, then is the corresponding *truncated normality* assumption equivalent to the original model?

As in the previous section, suppose that we have an exponential family \mathcal{F} with log-normalizer $g_{\mathcal{F}}$ over a support \mathcal{S} . For any distribution $D \in \mathcal{F}$, let $X \sim D$ be a random variable with density f_X and distribution function F_X .

Define another exponential family \mathcal{G} with the same sufficient statistics T and carrier measure h as \mathcal{F} , but over support $\mathcal{T} \subseteq \mathcal{S}$. Let Y be a random variable corresponding to X such that they have the same natural

parameters $\boldsymbol{\eta}$. Let its density be f_Y and its distribution function be F_Y . We have:

$$f_Y(y) = \frac{f_X(y)}{\mathcal{P}(X \in \mathcal{T})}. \quad (65)$$

The divisor $\mathcal{P}(X \in \mathcal{T})$ depends on the parameters $\boldsymbol{\eta}$, which means that \mathcal{G} has a different log-normalizer than \mathcal{F} :

$$g_{\mathcal{G}}(\boldsymbol{\eta}) = g_{\mathcal{F}}(\boldsymbol{\eta}) + \log(\mathcal{P}(X \in \mathcal{T})). \quad (66)$$

X and Y having different log-normalizers means that their expectation parameters $\boldsymbol{\chi}_X$ and $\boldsymbol{\chi}_Y$ are different even though their natural parameters are the same:

$$\boldsymbol{\chi}_Y \triangleq \mathbb{E}(T(Y)) \quad (\text{by equation 8}) \quad (67)$$

$$= \int_{\mathcal{T}} f_Y(y) T(y) dy. \quad (68)$$

This shows that clipping the distributional assumption of a generalized linear model changes the model.

1.5.3 Altering the carrier measure

Truncation (§1.5.2) of the sample space of an exponential family is equivalent to setting the carrier measure h to $-\infty$ over the truncated region. In the previous section, this would mean that we could have left \mathcal{G} 's support as \mathcal{S} , but set its carrier measure to

$$h_{\mathcal{G}}(x) = \begin{cases} h_{\mathcal{F}}(x) & \text{if } x \in \mathcal{T} \\ -\infty & \text{otherwise.} \end{cases} \quad (69)$$

If one desires a softer version of truncation, then one can softly decrease the carrier measure. This affects the expectation parameters—which are the expected value of the sufficient statistics—because it focuses that expectation where the carrier measure is larger. Therefore, altering the carrier measure changes the model.

2 Probability distributions

Listed below are many useful exponential families decomposed according to §1. See Nielsen and Garcia (2011) and Shao (2003, p. 66) for details.

A few liberties were taken with notation.

When the sufficient statistic is a matrix M , and the natural parameter is a matrix N , then the contribution to the log-probability is $\text{tr}(MN)$. When $M = \mathbf{x}\mathbf{y}^T$, $\text{tr}(MN) = \mathbf{x}^T N \mathbf{y}$.

When the sufficient statistic is complex $\mathbf{x} \in \mathbb{C}$ and/or the natural parameter is complex $\boldsymbol{\eta} \in \mathbb{C}$, then the contribution to the log-probability is $\text{Re}(\mathbf{x}^T \boldsymbol{\eta})$. This essentially means that there are

- two sufficient statistics $\text{Re}(\mathbf{x})$ and $\text{Im}(\mathbf{x})$, and
- two corresponding natural parameters $\text{Re}(\boldsymbol{\eta})$ and $-\text{Im}(\boldsymbol{\eta})$.

2.1 Normal distributions

The *normal distribution* arises as a result of the *central limit theorem*, which states that under mild conditions, the sum of many random variables will be approximately normally distributed. Kallenberg (2010) also characterizes this distribution using spherical symmetry in the context of *Gaussian processes*.

2.1.1 Univariate real

The univariate real normal distribution is the simplest normal distribution.

Table 1: The normal distribution.

density function	$f(x \mid \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$
sufficient statistic	$T(x) = (x, x^2)$
log-normalizer	$g(\boldsymbol{\eta}) = -\frac{\eta_1^2}{4\eta_2} + \frac{1}{2} \log\left(-\frac{\pi}{\eta_2}\right)$
carrier measure	$h(x) = 0$
support	$x \in \mathbb{R}$
Parameters	
source parameters	$(\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}_{\geq 0}$
natural parameters	$\boldsymbol{\eta} \in \mathbb{R} \times \mathbb{R}_{\leq 0}$
expectation parameters	$\boldsymbol{\chi} \in \mathbb{R} \times \mathbb{R}_{\geq 0}$
Parameter transformations	
source to natural parameters	$\boldsymbol{\eta} = \left(\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2}\right)$
source to expectation parameters	$\boldsymbol{\chi} = (\mu, \mu^2 + \sigma^2)$
natural to expectation parameters	$\boldsymbol{\chi} = \left(-\frac{\eta_1}{2\eta_2}, \frac{\eta_1^2}{4\eta_2^2} - \frac{1}{2\eta_2}\right)$

2.1.2 Real multivariate

The *multivariate normal distribution* is the generalization of the normal distribution (§2.1.1) to k dimensions.

We say that a vector distributed this way is *jointly normal*.

A few special cases are important exponential families:

- \mathbf{S} is fixed to \mathbf{I}_k ,
- \mathbf{S} is fixed to $\ell\mathbf{I}_k$ for some fixed ℓ ,
- \mathbf{S} equals $\ell\mathbf{I}_k$ for a parameter ℓ (isotropic variance), and
- \mathbf{S} equals $\text{diag } \ell$ for a parameter ℓ (diagonal variance).

Otherwise, we say that the distribution has general variance.

Table 2: The multivariate normal distribution.

density function	$f(\mathbf{x} \mid \boldsymbol{\mu}, \mathbf{S}) = \frac{1}{\sqrt{(2\pi)^k \det(\mathbf{S})}} \exp\left(-\frac{(\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{S}^{-1}(\mathbf{x} - \boldsymbol{\mu})}{2}\right)$
sufficient statistic	$T(\mathbf{x}) = (\mathbf{x}, \mathbf{x}\mathbf{x}^\top)$
log-normalizer	$g(\boldsymbol{\eta}, \mathbf{H}) = -\frac{\boldsymbol{\eta}^\top \mathbf{H}^{-1} \boldsymbol{\eta}}{4} - \frac{\log \det(-\mathbf{H})}{2} + \frac{k \log \pi}{2}$
carrier measure	$h(\mathbf{x}) = 0$
support	$\mathbf{x} \in \mathbb{R}^k$
Parameters	
source parameters ¹	$(\boldsymbol{\mu}, \mathbf{S}) \in \mathbb{R}^k \times \mathbb{R}_{\text{ps},s}^{k \times k}$
natural parameters ²	$(\boldsymbol{\eta}, \mathbf{H}) \in \mathbb{R}^k \times \mathbb{R}_{\text{ns},s}^{k \times k}$
expectation parameters ¹	$(\boldsymbol{\chi}, X) \in \mathbb{R}^k \times \mathbb{R}_{\text{ps},s}^{k \times k}$
Parameter transformations	
source to natural parameters	$(\boldsymbol{\eta}, \mathbf{H}) = (\mathbf{S}^{-1} \boldsymbol{\mu}, -\frac{1}{2} \mathbf{S}^{-1})$
source to expectation parameters	$(\boldsymbol{\chi}, X) = (\boldsymbol{\mu}, \boldsymbol{\mu} \boldsymbol{\mu}^\top + \mathbf{S})$
natural to expectation parameters	$(\boldsymbol{\chi}, X) = \left(-\frac{1}{2} \mathbf{H}^{-1} \boldsymbol{\eta}, \frac{1}{4} (\mathbf{H}^{-1} \boldsymbol{\eta}) (\mathbf{H}^{-1} \boldsymbol{\eta})^\top - \frac{1}{2} \mathbf{H}^{-1}\right)$

¹ $\mathbb{R}_{\text{ps},s}^{k \times k}$ is the set of positive semidefinite, symmetric $k \times k$ matrices of reals.

² $\mathbb{R}_{\text{ns},s}^{k \times k}$ is the set of negative semidefinite, symmetric $k \times k$ matrices of reals.

2.1.3 Complex multivariate

The *complex multivariate normal distribution* (Figure 4) is the generalization of the multivariate normal distribution (§2.1.2) to complex numbers (Picinbono 1996). This means that a vector of its real and imaginary components is *jointly normal* as per §2.1.2.

2.1.3.1 General

The general form of the complex multivariate normal distribution is given below.

One important special case has unit variance $\mathbf{S} = \mathbf{I}_k$ and zero pseudo-variance $\mathbf{U} = \mathbf{0}$.

Table 3: The complex multivariate normal distribution.

density function ⁴	$f(\mathbf{x} \mid \boldsymbol{\mu}, \mathbf{S}, \mathbf{U}) = \frac{\exp\left(-\frac{1}{2} \begin{bmatrix} \mathbf{x} - \boldsymbol{\mu} \\ \mathbf{x} - \boldsymbol{\mu} \end{bmatrix}^H \begin{bmatrix} \mathbf{S} & \mathbf{U} \\ \overline{\mathbf{U}} & \overline{\mathbf{S}} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{x} - \boldsymbol{\mu} \\ \mathbf{x} - \boldsymbol{\mu} \end{bmatrix}\right)}{\pi^k \sqrt{\det(\mathbf{S}) \det(\mathbf{P})}}$
sufficient statistic	$T(\mathbf{x}) = (\mathbf{x}, \mathbf{x}\mathbf{x}^H, \mathbf{x}\mathbf{x}^T)$
log-normalizer ^{6,7}	$g(\boldsymbol{\eta}, \mathbf{H}, \mathbf{J}) = -\boldsymbol{\mu}^H \mathbf{H} \boldsymbol{\mu} - \boldsymbol{\mu}^T \mathbf{J} \boldsymbol{\mu} + \frac{\log \det(\mathbf{S})}{2} - \frac{\log \det(-\mathbf{H})}{2} + k \log \pi$
carrier measure	$h(\mathbf{x}) = 0$
support	$\mathbf{x} \in \mathbb{C}^k$
Parameters	
source parameters ^{1,3}	$(\boldsymbol{\mu}, \mathbf{S}, \mathbf{U}) \in \mathbb{C}^k \times \mathbb{C}_{\text{ps,h}}^{k \times k} \times \mathbb{C}_s^{k \times k}$
natural parameters ^{2,3}	$(\boldsymbol{\eta}, \mathbf{H}, \mathbf{J}) \in \mathbb{C}^k \times \mathbb{C}_{\text{ns,h}}^{k \times k} \times \mathbb{C}_s^{k \times k}$
expectation parameters ^{1,3}	$(\boldsymbol{\chi}, \mathbf{X}, \mathbf{Y}) \in \mathbb{C}^k \times \mathbb{C}_{\text{ps,h}}^{k \times k} \times \mathbb{C}_s^{k \times k}$
Parameter transformations	
source to natural parameters ^{4,5}	$(\boldsymbol{\eta}, \mathbf{H}, \mathbf{J}) = \left(2 \left(\mathbf{P}^{-1} \bar{\boldsymbol{\mu}} - \mathbf{R}^T \overline{\mathbf{P}^{-1} \boldsymbol{\mu}}\right), -\overline{\mathbf{P}^{-1}}, \mathbf{R}^T \overline{\mathbf{P}^{-1}}\right)$
source to expectation parameters	$(\boldsymbol{\chi}, \mathbf{X}, \mathbf{Y}) = (\boldsymbol{\mu}, \boldsymbol{\mu} \boldsymbol{\mu}^H + \mathbf{S}, \boldsymbol{\mu} \boldsymbol{\mu}^T + \mathbf{U})$
natural to expectation parameters ^{6,7}	$(\boldsymbol{\chi}, \mathbf{X}, \mathbf{Y}) = (\boldsymbol{\mu}, \boldsymbol{\mu} \boldsymbol{\mu}^H + \mathbf{S}, \boldsymbol{\mu} \boldsymbol{\mu}^T + \mathbf{U})$

¹ $\mathbb{C}_{\text{ps,h}}^{k \times k}$ is the set of positive semidefinite, Hermitian $k \times k$ matrices of real numbers.

² $\mathbb{C}_{\text{ns,h}}^{k \times k}$ is the set of negative semidefinite, Hermitian $k \times k$ matrices of real numbers.

³ $\mathbb{C}_s^{k \times k}$ is the set of symmetric $k \times k$ matrices of complex numbers.

⁴ $\mathbf{P} = \overline{\mathbf{S}} - \mathbf{R} \mathbf{U} \in \mathbb{C}_{\text{ps,h}}^{k \times k}$

⁵ $\mathbf{R} = \mathbf{U}^H \mathbf{S}^{-1}$

⁶ $\mathbf{S} = (\overline{\mathbf{R}} \mathbf{R} - \mathbf{I}_k)^{-1} \mathbf{H}^{-1}$ and $\mathbf{U} = \overline{\mathbf{R}} \mathbf{S}$ where $\mathbf{R} = -(\mathbf{J} \mathbf{H}^{-1})^T$

⁷ $\boldsymbol{\mu} = \overline{\mathbf{L}} \boldsymbol{\eta} - (\mathbf{H}^{-1})^T \mathbf{J} \mathbf{L} \boldsymbol{\eta}$ where $\mathbf{L} = -\frac{(\mathbf{I}_k - \mathbf{K} \overline{\mathbf{K}})^{-1} (\mathbf{H}^{-1})^T}{2}$ and $\mathbf{K} = (\mathbf{H}^{-1})^T \mathbf{J}$

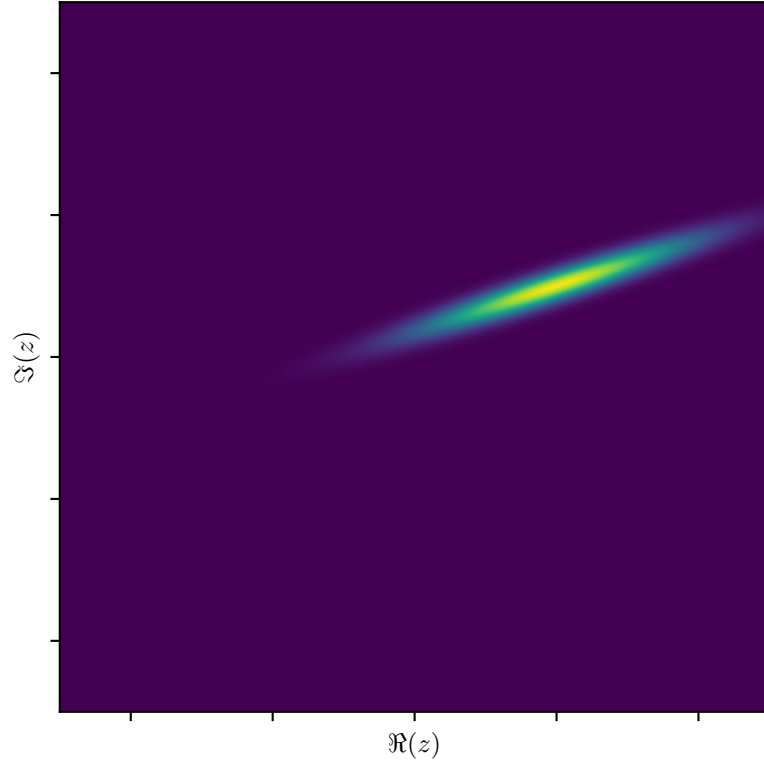


Figure 4: The density of a complex univariate normal distribution with a strong pseudo-curvature term that gives the distribution a definite phase.

2.1.3.2 Circularly-symmetric normal

The *circularly-symmetric normal distribution* is the special case of the complex multivariate normal distribution (§2.1.3) with an assumption of circular symmetry. Gallager (2008) proves that X is circularly-symmetric normally distributed if and only if it is jointly normally distributed, and

$$X = e^{i\theta} X \quad \forall \theta \in \mathbb{R}. \quad (70)$$

Table 4: The circularly-symmetric normal distribution.

density function	$f(\mathbf{x} \mid \mathbf{S}) = \frac{1}{\pi^k \det(\mathbf{S})} \exp(-\mathbf{x}^H \mathbf{S}^{-1} \mathbf{x})$
sufficient statistic	$T(\mathbf{x}) = \mathbf{x} \mathbf{x}^H$
log-normalizer	$g(\mathbf{H}) = -\log \det(-\mathbf{H}) + k \log \pi$
carrier measure	$h(\mathbf{x}) = 0$
support	$\mathbf{x} \in \mathbb{C}^k$
Parameters	
source parameters ¹	$\mathbf{S} \in \mathbb{C}_{\text{ps,h}}^{k \times k}$
natural parameters ²	$\mathbf{H} \in \mathbb{C}_{\text{ns,h}}^{k \times k}$
expectation parameters ¹	$\mathbf{X} \in \mathbb{C}_{\text{ps,h}}^{k \times k}$
Parameter transformations	
source to natural parameters	$\mathbf{H} = -\overline{\mathbf{S}^{-1}}$
source to expectation parameters	$\mathbf{X} = \mathbf{S}$
natural to expectation parameters	$\mathbf{X} = -\overline{\mathbf{H}^{-1}}$

¹ $\mathbb{C}_{\text{ps,h}}^{k \times k}$ is the set of positive semidefinite, Hermitian $k \times k$ matrices of complex numbers.

² $\mathbb{C}_{\text{ns,h}}^{k \times k}$ is the set of negative semidefinite, Hermitian $k \times k$ matrices of complex numbers.

2.2 Distributions on a finite set

2.2.1 Uniform (discrete)

The *discrete uniform distribution* arises when all points in its finite, discrete support Ω are equiprobable.

Table 5: The Uniform (discrete) distribution. (It has no parameters.)

density function	$f(\mathbf{x}) = \Omega ^{-1}$
sufficient statistic	$T(\mathbf{x}) = ()$
log-normalizer	$g(\boldsymbol{\eta}) = \log \Omega $
carrier measure	$h(\mathbf{x}) = 0$
support	$\mathbf{x} \in \Omega$

2.2.2 Multinomial, categorical, and Bernoulli

Suppose one draws n coloured balls from an urn (in which there are k different colours) replacing them between draws. The *multinomial distribution* is a probability distribution over the drawn colours.

The *categorical distribution* is the special case for $n = 1$; the *Binomial distribution* is the special case for $k = 2$; and the *Bernoulli distribution* is the special case for $n = 1$ and $k = 2$.

When n and k are fixed, the multinomial distribution is an exponential family.

Table 6: The multinomial distribution.

mass function	$f(\mathbf{x} \mid \mathbf{p}) = \frac{n!}{\prod_i x_i!} \prod_i p_i^{x_i}$
sufficient statistic	$T(\mathbf{x}) = (x_1, \dots, x_{k-1})$
log-normalizer	$g(\boldsymbol{\eta}) = \log \left(1 + \sum_{i=1}^{k-1} e^{\eta_i} \right) - \log n!$
carrier measure	$h(x) = -\sum_{i=1}^k \log x_i!$
support	$\mathbf{x} \in \{0, \dots, n\}^k$ with $\sum_i x_i = n$
Parameters	
source parameters	$\mathbf{p} \in [0, 1]^k$ where $\sum_i p_i = 1$
natural parameters	$\boldsymbol{\eta} \in \mathbb{R}^{k-1}$
expectation parameters	$\boldsymbol{\chi} \in [0, n]^{k-1}$
Parameter transformations	
source to natural parameters	$\boldsymbol{\eta} = (\log(p_i/p_k))_{i=1}^{k-1}$
source to expectation parameters	$\boldsymbol{\chi} = (np_i)_{i=1}^{k-1}$
natural to expectation parameters ¹	$\boldsymbol{\chi} = \left(\frac{ne^{\eta_i}}{A} \right)_{i=1}^{k-1}$

¹ where

$$A \triangleq 1 + \sum_{i=1}^{k-1} e^{\eta_i} = \sum_{i=1}^k \frac{p_i}{p_k}$$

2.3 Distributions on the nonnegative integers

2.3.1 Negative binomial and geometric

The *negative binomial distribution* models the number of failures before r successes of a Bernoulli distribution having probability p (§2.2.2). The *geometric distribution* is the special case when $r = 1$. As long as r is fixed,

the negative binomial distribution is an exponential family.

Table 7: The negative binomial distribution.

mass function	$f(x p) = \binom{x+r-1}{x} (1-p)^x p^r$
sufficient statistic	$T(x) = x$
log-normalizer	$g(\eta) = -r \log(1 - e^\eta)$
carrier measure	$h(x) = \log \binom{x+r-1}{x}$
support	$\mathbf{x} \in \mathbb{Z}_{\geq 0}$
Parameters	
source parameters	$p \in [0, 1]$
natural parameters	$\eta \in \mathbb{R}_{\leq 0}$
expectation parameters	$\chi \in \mathbb{R}_{\geq 0}$
Parameter transformations	
source to natural parameters	$\eta = \log(1 - p)$
source to expectation parameters	$\chi = \left(\frac{1-p}{p}\right) r$
natural to expectation parameters	$\chi = \frac{r}{e^{-\eta} - 1}$

2.3.2 Logarithmic

A Poisson distribution compounded with a *logarithmic* distribution yields a negative binomial distribution.

Table 8: The logarithmic distribution.

mass function	$f(x p) = \frac{-p^x}{x \log(1 - p)}$
sufficient statistic	$T(x) = x$
log-normalizer	$g(\eta) = \log(-\log(1 - e^\eta))$
carrier measure	$h(x) = -\log x$
support	$\mathbf{x} \in \mathbb{Z}_{\geq 1}$
Parameters	
source parameters	$p \in [0, 1]$
natural parameters	$\eta \in \mathbb{R}_{\leq 0}$
expectation parameters	$\chi \in \mathbb{R}_{\geq 1}$
Parameter transformations	
source to natural parameters	$\eta = \log p$
source to expectation parameters	$\chi = \frac{-p}{(1-p) \log(1-p)}$
natural to expectation parameters	$\chi = \frac{-e^\eta}{(1-e^\eta) \log(1-e^\eta)}$

2.3.3 Poisson

The *Poisson distribution* arises in Poisson processes as described in §??.

Table 9: The Poisson distribution.

mass function	$f(x \mid \lambda) = \frac{\lambda^x e^{-\lambda}}{x!}$
sufficient statistic	$T(x) = x$
log-normalizer	$g(\eta) = \exp \eta$
carrier measure	$h(x) = -\log(x!)$
support	$\mathbf{x} \in \mathbb{Z}_{\geq 0}$
Parameters	
source parameters	$\lambda \in \mathbb{R}_{\geq 0}$
natural parameters	$\eta \in \mathbb{R}$
expectation parameters	$\chi \in \mathbb{R}_{\geq 0}$
Parameter transformations	
source to natural parameters	$\eta = \log \lambda$
source to expectation parameters	$\chi = \lambda$
natural to expectation parameters	$\chi = \exp \eta$

2.4 Distributions on the positive reals

2.4.1 Chi, chi-Square, exponential, gamma, Rayleigh, and Weibull

The *exponential distribution* with rate λ arises as the waiting time until the next occurrence of a linear Poisson process with intensity λ (§??). The *gamma distribution* is the waiting time for k occurrences, i.e., it is the sum of k exponentially-distributed random variables each having rate λ .

The inverse-gamma, chi-square, inverse-chi-square, chi, Weibull and Rayleigh distributions are all transformations or special cases of the gamma distribution:

If X is gamma-distributed with parameters k, λ , then:

- X^{-1} is inverse-gamma distributed with the same parameters.

If X is gamma-distributed with parameters $k, \lambda = 1/2$, then:

- X is chi-square distributed with parameter $2k$, and
- X^{-1} is inverse-chi-square distributed with parameter $2k$, and
- \sqrt{X} is chi distributed with parameter $2k$.

If X is exponentially-distributed with rate λ , then:

- \sqrt{X} is Rayleigh-distributed with parameter $\sqrt{\frac{\lambda}{2}}$, and

If X is exponentially-distributed with rate 1, then:

- $\sqrt[k]{X}$ is Weibull-distributed with shape k .

Table 10: The gamma distribution.

density function	$f(x k, \lambda) = \frac{\lambda^k}{\Gamma(k)} x^{k-1} e^{-\lambda x}$
sufficient statistic	$T(x) = (x, \log x)$
log-normalizer	$g(\boldsymbol{\eta}) = \log \Gamma(\eta_2 + 1) - (\eta_2 + 1) \log(-\eta_1)$
carrier measure	$h(x) = 0$
support	$\mathbf{x} \in \mathbb{R}_{\geq 0}$
Parameters	
source parameters	$(k, \lambda) \in \mathbb{R}_{\geq 0} \times \mathbb{R}_{\geq 0}$
natural parameters	$\boldsymbol{\eta} \in \mathbb{R}_{\leq 0} \times \mathbb{R}_{\geq -1}$
expectation parameters	$\boldsymbol{\chi} \in \mathbb{R}_{\geq 0} \times \mathbb{R}$
Parameter transformations	
source to natural parameters	$\boldsymbol{\eta} = (-\lambda, k - 1)$
source to expectation parameters	$\boldsymbol{\chi} = (k/\lambda, \psi(k) - \log(\lambda))$
natural to expectation parameters	$\boldsymbol{\chi} = \left(-\frac{\eta_2 + 1}{\eta_1}, \psi(\eta_2 + 1) - \log(-\eta_1) \right)$

¹ ψ is the digamma function: $\frac{d \log \Gamma(x)}{dx}$

2.5 Distributions on the simplex

2.5.1 Dirichlet, beta, and continuous uniform

Suppose one draws n coloured balls from an urn (in which there are k different colours) replacing them between draws. The *Dirichlet distribution* is often used to model the belief about the proportions of the coloured balls in the urn because it is the conjugate prior distribution (§1.3.1) of the multinomial distribution (§2.2.2).

The *beta distribution* is the special case of the Dirichlet distribution for $k = 2$.

The *continuous uniform distribution* arises when all points in its finite support $[0, 1]$ are equiprobable. It

is the special case of the beta distribution for $\boldsymbol{\alpha} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$.

Table 11: The Dirichlet distribution.

density function ¹	$f(\mathbf{x} \mid \boldsymbol{\alpha}) = \frac{1}{B(\boldsymbol{\alpha})} \prod_i x_i^{\alpha_i - 1}$
sufficient statistic	$T(\mathbf{x}) = (\log x_i)_i$
log-normalizer	$g(\boldsymbol{\eta}) = \sum_i \log \Gamma(\eta_i + 1) - \log \Gamma(\sum_j \eta_j + k)$
carrier measure	$h(x) = 0$
support	$\mathbf{x} \in [0, 1]^k$ with $\sum_i x_i = 1$
Parameters	
source parameters	$\boldsymbol{\alpha} \in \mathbb{R}_{\geq 0}^k$
natural parameters	$\boldsymbol{\eta} \in \mathbb{R}_{\geq -1}^k$
expectation parameters	$\boldsymbol{\chi} \in \mathbb{R}_{\leq 0}^k$
Parameter transformations	
source to natural parameters	$\boldsymbol{\eta} = (\alpha_i - 1)_i$
source to expectation parameters ²	$\boldsymbol{\chi} = (\psi(\alpha_i) - \psi(\sum_j \alpha_j))_i$
natural to expectation parameters ²	$\boldsymbol{\chi} = (\psi(\eta_i + 1) - \psi(\sum_j \eta_j + n))_i$

¹ B is the multinomial beta function: $B(\boldsymbol{\alpha}) = \frac{\prod_i \Gamma(\alpha_i)}{\Gamma(\sum_j \alpha_j)}$

² ψ is the digamma function: $\frac{d \log \Gamma(x)}{dx}$

2.5.2 Generalized Dirichlet

The *generalized Dirichlet distribution* (Wong 1998) has a more general covariance structure than the Dirichlet distribution.

Table 12: The generalized Dirichlet distribution.

density function ¹	$f(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \prod_{i=1}^k \frac{x_i^{\alpha_i-1} \left(1 - \sum_{j \leq i} x_j\right)^{\gamma_i}}{B(\alpha_i, \beta_i)}$
sufficient statistic	$T(\mathbf{x}) = (\log x_i)_i, (1 - \sum_{j \leq i} x_j)_i$
log-normalizer	$g(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \sum_i \log B(\alpha_i, \beta_i)$
carrier measure	$h(x) = 0$
support	$\mathbf{x} \in [0, 1]^k \text{ with } \sum_i x_i = 1$
Parameters	
source parameters	$(\boldsymbol{\alpha}, \boldsymbol{\beta}) \in \mathbb{R}_{\geq 0}^k \times \mathbb{R}_{\geq 0}^k$
natural parameters	$(\boldsymbol{\alpha}', \boldsymbol{\gamma}) \in \mathbb{R}_{\geq -1}^k \times \mathbb{R}^k$
expectation parameters	$(\boldsymbol{\chi}, \boldsymbol{\psi}) \in \mathbb{R}_{\leq 0}^k \times \mathbb{R}_{\leq 0}^k$
Parameter transformations	
source to natural parameters	$\alpha'_i = \alpha_i - 1$
	$\gamma_i = \begin{cases} \beta_i - \alpha_{i+1} - \beta_{i+1} & \text{if } i < k \\ \beta_k - 1 & \text{otherwise.} \end{cases}$
source to expectation parameters ²	$\chi_i = \alpha_i + \sum_{j < i} b_j$
	$\psi_i = \sum_{j \leq i} b_j$
natural to source parameters	$\alpha_i = \alpha'_i + 1$
	$\beta_i = \sum_{j \geq i} \gamma_j + \sum_{j > i} \alpha_j + 1$

¹ B is the multinomial beta function: $B(\boldsymbol{\alpha}) = \frac{\prod_i \Gamma(\alpha_i)}{\Gamma\left(\sum_j \alpha_j\right)}$

² Where ψ is the digamma function: $\frac{d \log \Gamma(x)}{dx}$ and

$$s_i = \psi(\alpha_i + \beta_i)$$

$$a_i = \psi(\alpha_i) - s_i$$

$$b_i = \psi(\beta_i) - s_i$$

2.6 Distributions on the n-sphere

2.6.1 Von Mises-Fisher and von Mises

The *von Mises-Fisher* distribution is a probability distribution on the k -dimensional unit sphere (Dhillon and Sra (2003)). When $k = 2$, the distribution collapses to the *von Mises* distribution on the unit circle.

Table 13: The von Mises-Fisher distribution.

density function ²	$f(\mathbf{x} \mid \kappa, \boldsymbol{\mu}) = c_k(\kappa) e^{\kappa \boldsymbol{\mu}^\top \mathbf{x}}$
sufficient statistic	$T(\mathbf{x}) = \mathbf{x}$
log-normalizer	$g(\boldsymbol{\eta}) = -\log c_k(\ \boldsymbol{\eta}\)$
carrier measure	$h(x) = 0$
support	$\mathbf{x} \in \mathbb{R}^k$ with $\ \mathbf{x}_i\ = 1$
Parameters	
source parameters	$\kappa, \boldsymbol{\mu} \in \mathbb{R}_{\geq 0} \times \mathbb{R}^k$ with $\ \boldsymbol{\mu}\ = 1$
natural parameters	$\boldsymbol{\eta} \in \mathbb{R}^k$
expectation parameters	$\boldsymbol{\chi} \in \mathbb{R}^k$ with $\ \boldsymbol{\chi}\ \leq 1$
Parameter transformations	
source to natural parameters	$\boldsymbol{\eta} = \kappa \boldsymbol{\mu}$
source to expectation parameters ³	$\boldsymbol{\chi} = A_k(\kappa) \boldsymbol{\mu}$
natural to expectation parameters	$\boldsymbol{\chi} = A_k(\ \boldsymbol{\eta}\) \frac{\boldsymbol{\eta}}{\ \boldsymbol{\eta}\ }$

¹ I_k is a modified Bessel function of the first kind at order k :

$$I_k(\kappa) = \sum_{i \geq 0} \frac{\left(\frac{\kappa}{2}\right)^{2i+k}}{\Gamma(i+k+1)i!}.$$

² The reciprocal normalizer is

$$c_k(\kappa) = \frac{\kappa^{k/2-1}}{(2\pi)^{k/2} I_{k/2-1}(\kappa)}.$$

³ A is defined

$$A_k(\kappa) = \frac{I_{k/2}(\kappa)}{I_{k/2-1}(\kappa)}$$

and

$$A_k^{-1}(\mu) \approx \frac{\mu k - \mu^3}{1 - \mu^2}.$$

References

- [1] I. S. Dhillon and S. Sra, “Modeling data using directional distributions,” University of Texas, 2003, pp. 1–21.
- [2] B. Efron, “The geometry of exponential families,” *The Annals of Statistics*, vol. 6, pp. 362–376, 2 1978.
- [3] R. G. Gallager, “Circularly-symmetric gaussian random vectors,” *preprint*, pp. 1–9, 2008.

- [4] D. B. Gokhale, *Maximum entropy characterizations of some distributions*, G. P. Patil, S. Kotz, and J. K. Ord, Eds., 1975.
- [5] G. E. Hinton, “Training products of experts by minimizing contrastive divergence,” *Neural Computation*, vol. 14, pp. 1771–1800, 8 Aug. 2002.
- [6] E. T. Jaynes, “Information theory and statistical mechanics,” *Phys. Rev.*, vol. 106, pp. 620–630, 4 May 1957.
- [7] O. Kallenberg, *Foundations of Modern Probability*. Springer, 2010.
- [8] F. Nielsen and V. Garcia, “Statistical exponential families: A digest with flash cards,” *CoRR*, vol. abs/0911.4, 2011.
- [9] F. Nielsen and R. Nock, “Entropies and cross-entropies of exponential families,” *Proceedings of the International Conference on Image Processing, ICIP 2010, September 26-29, Hong Kong, China*, pp. 3621–3624, 2011.
- [10] B. Picinbono, “Second-order complex random vectors and normal distributions,” *IEEE Transactions on Signal Processing*, vol. 44, pp. 2637–2640, 10 1996.
- [11] J. Shao, *Mathematical Statistics*. Springer, 2003.
- [12] T.-T. Wong, “Generalized dirichlet distribution in bayesian analysis,” 1998, pp. 165–181.